



# Towards Static Analysis of Functional Programs using Tree Automata Completion

Thomas Genet

## ► To cite this version:

Thomas Genet. Towards Static Analysis of Functional Programs using Tree Automata Completion.  
[Research Report] 2013, pp.15. hal-00921814

**HAL Id: hal-00921814**

**<https://inria.hal.science/hal-00921814>**

Submitted on 21 Dec 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Towards Static Analysis of Functional Programs using Tree Automata Completion

Thomas Genet

INRIA/IRISA, Université de Rennes, France  
genet@irisa.fr

**Abstract.** Tree Automata Completion is a family of techniques for computing or approximating the set of terms reachable by a rewriting relation. The completion algorithm we focus on is parameterized by a set  $E$  of equations controlling the precision of the approximation and influencing its termination. In this work, we study the application of completion to the static analysis of functional programs. We give a sufficient condition on  $\mathcal{T}(\mathcal{F})$  and  $E$  for completion algorithm to always terminate. In the particular setting of functional programs, this condition can be relaxed into a condition on  $\mathcal{T}(\mathcal{C})$  and  $E$  that is close to what is generally done in static analysis where abstractions are performed on data. For functional programs translated into TRSs, we give a sufficient condition for completion to terminate.

## 1 Introduction

This work presents a first step towards the application of reachability analysis techniques coming from rewriting to the static analysis of functional programs. Computing or approximating the set of terms reachable by rewriting has more and more applications. For a Term Rewriting System (TRS)  $\mathcal{R}$  and a set of terms  $L_0 \subseteq \mathcal{T}(\mathcal{F})$ , the set of reachable terms is  $R^*(L_0) = \{t \in \mathcal{T}(\mathcal{F}) \mid \exists s \in L_0, s \rightarrow_{\mathcal{R}}^* t\}$ . This set can be computed exactly for specific classes of  $\mathcal{R}$  [12] but, in general, it has to be approximated. Applications of the approximation of  $R^*(L_0)$  are ranging from cryptographic protocol verification [13, 1], to static analysis of various programming languages [5, 22] or to TRS termination proofs [24, 18]. Most of the techniques compute such approximations using tree automata as the core formalism to represent or approximate the (possibly) infinite set of terms  $R^*(L_0)$ . Most of them also rely on a Knuth-Bendix completion-like algorithm completing a tree automaton  $\mathcal{A}$  recognizing  $L_0$  into an automaton  $\mathcal{A}^*$  recognizing exactly, or over-approximating, the set  $R^*(L_0)$ . As a result, these techniques can be referred as *tree automata completion* techniques [11, 27, 26, 10, 4, 15, 23]. A strength of this algorithm, and at the same time a weakness, is that its precision is parameterized by a function [10] or a set of equations [15]. It is a strength because tuning the approximation function (or equations) permits to adapt the precision of completion to a specific goal to tackle. This is what made it successful for program and protocol verification. On the other hand, this is a weakness because it is difficult to guarantee its termination.

In this paper, we define a simple sufficient condition on the set of equations for the tree automata completion algorithm to terminate. This condition, which is strong in general, reveals to be natural and well adapted for the approximation of TRSs encoding typed functional programs. We thus obtain a way to automatically over-approximate the set of all reachable program states of a functional program, or even restrict it to the set of all results. Thus we can over-approximate the co-domain of a functional program.

## 2 Related work

*Tree automata completion.* With regards to most papers about completion [11, 27, 26, 10, 4, 15, 23], our contribution is to give the first criterion *on the approximation* for the completion to terminate. Note that it is possible to guarantee termination of the completion by inferring an approximation

adapted to the TRS under concern, like in [6]. In this case, given a TRS, the approximation is fixed and unique. Our solution is more flexible because it lets the user change the precision of the approximation while keeping the termination guarantee. In [26], T. Takai have a completion parameterized by a set of equations. He also gives a termination proof for its completion but only for some restricted classes of TRSs. Here our termination proof holds for any left-linear TRS provided that the set of equations satisfy some properties.

*Static analysis of functional programs.* With regards to static analysis of functional programs using grammars or automata, our contribution is in the scope of data-flow analysis techniques, rather than control-flow analysis. In other words, we are interested here in predicting the results of a function [25], rather than predicting the control flow [21]. Those two papers, as well as many other ones, deal with higher order functions using complex higher-order grammar formalisms (PMRS and HORS). Higher-order functions are not in the scope of the solution we propose here. However, we obtained some preliminary results suggesting that an extension to higher order functions is possible and gives relevant results (see Section 6). Furthermore, using equations, approximations are defined in a more declarative and flexible way than in [25], where they are defined by a dedicated algorithm. Besides, the verification mechanisms of [25] use automatic abstraction refinement. This can be also performed in the completion setting [3] and adapted to the analysis of functional programs [17]. Finally, using a simpler (first order) formalism, *i.e.* tree automata, makes it easier to take into account some other aspects like: evaluation strategies and built-ins types (see Section 6) that are not considered by those papers.

### 3 Background

In this section, we introduce some definitions and concepts that will be used throughout the rest of the paper (see also [2, 9]). Let  $\mathcal{F}$  be a finite set of symbols, each associated with an arity function, and let  $\mathcal{X}$  be a countable set of *variables*.  $\mathcal{T}(\mathcal{F}, \mathcal{X})$  denotes the set of *terms* and  $\mathcal{T}(\mathcal{F})$  denotes the set of *ground terms* (terms without variables). The set of variables of a term  $t$  is denoted by  $\text{Var}(t)$ . A *substitution* is a function  $\sigma$  from  $\mathcal{X}$  into  $\mathcal{T}(\mathcal{F}, \mathcal{X})$ , which can be uniquely extended to an endomorphism of  $\mathcal{T}(\mathcal{F}, \mathcal{X})$ . A *position*  $p$  for a term  $t$  is a word over  $\mathbb{N}$ . The empty sequence  $\lambda$  denotes the top-most position. The set  $\text{Pos}(t)$  of positions of a term  $t$  is inductively defined by  $\text{Pos}(t) = \{\lambda\}$  if  $t \in \mathcal{X}$  and  $\text{Pos}(f(t_1, \dots, t_n)) = \{\lambda\} \cup \{i.p \mid 1 \leq i \leq n \text{ and } p \in \text{Pos}(t_i)\}$  otherwise. If  $p \in \text{Pos}(t)$ , then  $t|_p$  denotes the subterm of  $t$  at position  $p$  and  $t[s]_p$  denotes the term obtained by replacement of the subterm  $t|_p$  at position  $p$  by the term  $s$ .

A *term rewriting system* (TRS)  $\mathcal{R}$  is a set of *rewrite rules*  $l \rightarrow r$ , where  $l, r \in \mathcal{T}(\mathcal{F}, \mathcal{X})$ ,  $l \notin \mathcal{X}$ , and  $\text{Var}(l) \supseteq \text{Var}(r)$ . A rewrite rule  $l \rightarrow r$  is *left-linear* if each variable of  $l$  occurs only once in  $l$ . A TRS  $\mathcal{R}$  is left-linear if every rewrite rule  $l \rightarrow r$  of  $\mathcal{R}$  is left-linear. The TRS  $\mathcal{R}$  induces a rewriting relation  $\rightarrow_{\mathcal{R}}$  on terms as follows. Let  $s, t \in \mathcal{T}(\mathcal{F}, \mathcal{X})$  and  $l \rightarrow r \in \mathcal{R}$ ,  $s \rightarrow_{\mathcal{R}} t$  denotes that there exists a position  $p \in \text{Pos}(s)$  and a substitution  $\sigma$  such that  $s|_p = l\sigma$  and  $t = s[r\sigma]_p$ . Given a TRS  $\mathcal{R}$ , there is a partition  $(\mathcal{C}, \mathcal{D})$  of  $\mathcal{F}$  such that all symbols occurring at the root position of left-hand sides of rules of  $\mathcal{R}$  are in  $\mathcal{D}$ .  $\mathcal{D}$  is the set of defined symbols of  $\mathcal{R}$ ,  $\mathcal{C}$  is the set of constructors. Terms in  $\mathcal{T}(\mathcal{C})$  are called data-terms. The reflexive transitive closure of  $\rightarrow_{\mathcal{R}}$  is denoted by  $\rightarrow_{\mathcal{R}}^*$  and  $s \rightarrow_{\mathcal{R}}^! t$  denotes that  $s \rightarrow_{\mathcal{R}}^* t$  and  $t$  is irreducible by  $\mathcal{R}$ . The set of irreducible terms w.r.t. a TRS  $\mathcal{R}$  is denoted by  $\text{IRR}(\mathcal{R})$ . The set of  $\mathcal{R}$ -descendants of a set of ground terms  $I$  is  $\mathcal{R}^*(I) = \{t \in \mathcal{T}(\mathcal{F}) \mid \exists s \in I \text{ s.t. } s \rightarrow_{\mathcal{R}}^* t\}$ . A TRS  $\mathcal{R}$  is sufficiently complete if for all  $s \in \mathcal{T}(\mathcal{F})$ ,  $(\mathcal{R}^*(\{s\}) \cap \mathcal{T}(\mathcal{C})) \neq \emptyset$ .

An *equation set*  $E$  is a set of *equations*  $l = r$ , where  $l, r \in \mathcal{T}(\mathcal{F}, \mathcal{X})$ . The relation  $=_E$  is the smallest congruence such that for all substitution  $\sigma$  we have  $l\sigma = r\sigma$ . Given a TRS  $\mathcal{R}$  and a set of equations  $E$ , a term  $s \in \mathcal{T}(\mathcal{F})$  is rewritten modulo  $E$  into  $t \in \mathcal{T}(\mathcal{F})$ , denoted  $s \rightarrow_{\mathcal{R}/E} t$ , if there exist  $s' \in \mathcal{T}(\mathcal{F})$  and  $t' \in \mathcal{T}(\mathcal{F})$  such that  $s =_E s' \rightarrow_{\mathcal{R}} t' =_E t$ . The reflexive transitive closure  $\rightarrow_{\mathcal{R}/E}^*$  of  $\rightarrow_{\mathcal{R}/E}$  is defined as usual except that reflexivity is extended to terms equal modulo  $E$ ,

i.e. for all  $s, t \in \mathcal{T}(\mathcal{F})$  if  $s =_E t$  then  $s \rightarrow_{\mathcal{R}/E}^* t$ . The set of  $\mathcal{R}$ -descendants modulo  $E$  of a set of ground terms  $I$  is  $\mathcal{R}_E^*(I) = \{t \in \mathcal{T}(\mathcal{F}) \mid \exists s \in I \text{ s.t. } s \rightarrow_{\mathcal{R}/E}^* t\}$ .

Let  $\mathcal{Q}$  be a countably infinite set of symbols with arity 0, called *states*, such that  $\mathcal{Q} \cap \mathcal{F} = \emptyset$ .  $\mathcal{T}(\mathcal{F} \cup \mathcal{Q})$  is called the set of *configurations*. A *transition* is a rewrite rule  $c \rightarrow q$ , where  $c$  is a configuration and  $q$  is state. A transition is *normalized* when  $c = f(q_1, \dots, q_n)$ ,  $f \in \mathcal{F}$  is of arity  $n$ , and  $q_1, \dots, q_n \in \mathcal{Q}$ . An  $\epsilon$ -transition is a transition of the form  $q \rightarrow q'$  where  $q$  and  $q'$  are states. A bottom-up non-deterministic finite tree automaton (tree automaton for short) over the alphabet  $\mathcal{F}$  is a tuple  $\mathcal{A} = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_F, \Delta \rangle$ , where  $\mathcal{Q}_F \subseteq \mathcal{Q}$ ,  $\Delta$  is a set of normalized transitions and  $\epsilon$ -transitions. The transitive and reflexive *rewriting relation* on  $\mathcal{T}(\mathcal{F} \cup \mathcal{Q})$  induced by the set of transitions  $\Delta$  (resp. all transitions except  $\epsilon$ -transitions) is denoted by  $\rightarrow_\Delta^*$  (resp.  $\rightarrow_\Delta^{\epsilon*}$ ). When  $\Delta$  is attached to a tree automaton  $\mathcal{A}$  we also note those two relations  $\rightarrow_{\mathcal{A}}^*$  and  $\rightarrow_{\mathcal{A}}^{\epsilon*}$ , respectively. A tree automaton  $\mathcal{A}$  is complete if for all  $s \in \mathcal{T}(\mathcal{F})$  there exists a state  $q$  of  $\mathcal{A}$  such that  $s \rightarrow_{\mathcal{A}}^* q$ . The language (resp.  $\epsilon$ -language) recognized by  $\mathcal{A}$  in a state  $q$  is  $\mathcal{L}(\mathcal{A}, q) = \{t \in \mathcal{T}(\mathcal{F}) \mid t \rightarrow_{\mathcal{A}}^* q\}$  (resp.  $\mathcal{L}^{\epsilon}(\mathcal{A}, q) = \{t \in \mathcal{T}(\mathcal{F}) \mid t \rightarrow_{\mathcal{A}}^{\epsilon*} q\}$ ). A state  $q$  of an automaton  $\mathcal{A}$  is *reachable* (resp.  $\epsilon$ -reachable) if  $\mathcal{L}(\mathcal{A}, q) \neq \emptyset$  (resp.  $\mathcal{L}^{\epsilon}(\mathcal{A}, q) \neq \emptyset$ ). We define  $\mathcal{L}(\mathcal{A}) = \bigcup_{q \in \mathcal{Q}_F} \mathcal{L}(\mathcal{A}, q)$ . A set of transitions  $\Delta$  is  $\epsilon$ -deterministic there are no two normalized transitions in  $\Delta$  with the same left-hand side. A tree automaton  $\mathcal{A}$  is  $\epsilon$ -deterministic if its set of transition is  $\epsilon$ -deterministic. Note that if  $\mathcal{A}$  is  $\epsilon$ -deterministic then for all states  $q_1, q_2$  of  $\mathcal{A}$  such that  $q_1 \neq q_2$ , we have  $\mathcal{L}^{\epsilon}(\mathcal{A}, q_1) \cap \mathcal{L}^{\epsilon}(\mathcal{A}, q_2) = \emptyset$ .

## 4 Tree Automata Completion Algorithm

Tree Automata Completion algorithms were proposed in [19, 11, 27, 15]. Tree automata completion is very similar to a Knuth-Bendix completion except that it runs on two distinct sets of rules: a TRS  $\mathcal{R}$  and a set of transitions  $\Delta$  of a tree automaton  $\mathcal{A}$ . A critical pair is defined as follows: if there is a rewrite rule  $l \rightarrow r \in \mathcal{R}$  and a substitution  $\sigma : \mathcal{X} \mapsto \mathcal{Q}$  such that

$$\begin{array}{ccc} l\sigma & \xrightarrow{\mathcal{R}} & r\sigma \\ \Delta \downarrow^* & & \\ q & & \end{array}$$

then we know that  $l\sigma$  is rewritten (recognized) by  $\mathcal{A}$  (rules of  $\Delta$ ) and that  $l\sigma$  is rewritten into  $r\sigma$  by  $\mathcal{R}$ . If  $r\sigma \not\rightarrow_\Delta^* q$  then the critical pair has to be solved for  $r\sigma$  to be recognized by  $\Delta$  into state  $q$ . Hence, we need to add the necessary transitions to  $\Delta$  to have  $r\sigma \rightarrow_\Delta^* q$ . Note that, contrary to Knuth-Bendix completion, we do not have any choice here w.r.t. the direction since having  $q \rightarrow_\Delta^* r\sigma$  is not compatible with the standard transition relation of a tree automata. Then, as in Knuth-Bendix completion, this process is iterated until all critical pairs between  $\mathcal{R}$  and  $\Delta$  can be joined. The complete process is described in the following section.

### 4.1 Tree Automata Completion General Principle

Let us first recall the tree automata completion principle. Starting from a tree automaton  $\mathcal{A}_0 = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_F, \Delta_0 \rangle$  and a left-linear TRS  $\mathcal{R}$ , the aim of the completion algorithm is to compute a tree automaton  $\mathcal{A}'$  such that  $\mathcal{L}(\mathcal{A}') = \mathcal{R}^*(\mathcal{L}(\mathcal{A}_0))$  or  $\mathcal{L}(\mathcal{A}') \supseteq \mathcal{R}^*(\mathcal{L}(\mathcal{A}_0))$ . Then  $\mathcal{A}'$  is used to show that terms recognized by a tree automaton  $\mathcal{A}_{bad}$  are not reachable by rewriting terms of  $\mathcal{L}(\mathcal{A}_0)$  with  $\mathcal{R}$ , i.e.  $\forall s \in \mathcal{L}(\mathcal{A}_0) \forall t \in \mathcal{L}(\mathcal{A}_{bad}) : s \not\rightarrow_{\mathcal{R}}^* t$ . For this, it is enough to show that  $\mathcal{L}(\mathcal{A}') \cap \mathcal{L}(\mathcal{A}_{bad}) = \emptyset$ , i.e. compute the automaton recognizing the intersection and show that the recognized language is empty.

Tree automata completion successively computes tree automata  $\mathcal{A}_{\mathcal{R}}^1, \mathcal{A}_{\mathcal{R}}^2, \dots$  such that  $\forall i \geq 0 : \mathcal{L}(\mathcal{A}_{\mathcal{R}}^i) \subseteq \mathcal{L}(\mathcal{A}_{\mathcal{R}}^{i+1})$  and if  $s \in \mathcal{L}(\mathcal{A}_{\mathcal{R}}^i)$ , such that  $s \rightarrow_{\mathcal{R}} t$  then  $t \in \mathcal{L}(\mathcal{A}_{\mathcal{R}}^{i+1})$ , until we get an

automaton  $\mathcal{A}_{\mathcal{R}}^k$  with  $k \in \mathbb{N}$  such that  $\mathcal{L}(\mathcal{A}_{\mathcal{R}}^k) = \mathcal{L}(\mathcal{A}_{\mathcal{R}}^{k+1})$ . Thus,  $\mathcal{A}_{\mathcal{R}}^k$  is a fixpoint and  $\mathcal{A}_{\mathcal{R}}^k$  also verifies  $\mathcal{L}(\mathcal{A}_{\mathcal{R}}^k) \supseteq \mathcal{R}^*(\mathcal{L}(\mathcal{A}_0))$ . To construct  $\mathcal{A}_{\mathcal{R}}^{i+1}$  from  $\mathcal{A}_{\mathcal{R}}^i$ , we achieve a *completion step* which consists in finding *critical pairs* between  $\rightarrow_{\mathcal{R}}$  and  $\rightarrow_{\mathcal{A}_{\mathcal{R}}^i}$ . For a substitution  $\sigma : \mathcal{X} \mapsto \mathcal{Q}$  and a rule  $l \rightarrow r \in \mathcal{R}$ , a critical pair is an instance  $l\sigma$  of  $l$  such that there exists  $q \in \mathcal{Q}$  satisfying  $l\sigma \rightarrow_{\mathcal{A}_{\mathcal{R}}^i}^* q$  and  $r\sigma \not\rightarrow_{\mathcal{A}_{\mathcal{R}}^i}^* q$ . For  $r\sigma$  to be recognized by the same state and thus model the rewriting of  $l\sigma$  into  $r\sigma$ , it is enough to add the necessary transitions to  $\mathcal{A}_{\mathcal{R}}^i$  to obtain  $\mathcal{A}_{\mathcal{R}}^{i+1}$  such that  $r\sigma \rightarrow_{\mathcal{A}_{\mathcal{R}}^{i+1}}^* q$ . In [27, 15], critical pairs are joined in the following way:

$$\begin{array}{ccc} l\sigma & \xrightarrow{\mathcal{R}} & r\sigma \\ \mathcal{A}_{\mathcal{R}}^i \downarrow & & \downarrow \mathcal{A}_{\mathcal{R}}^{i+1} \\ q & \xleftarrow{\mathcal{A}_{\mathcal{R}}^{i+1}} & q' \end{array}$$

From an algorithmic point of view, there remains two interesting points: how to find the  $\sigma$  substitutions, *i.e.* how to perform matching? and how to find the necessary transitions to add to  $\mathcal{A}_{\mathcal{R}}^i$  in order to have  $r\sigma \rightarrow_{\mathcal{A}_{\mathcal{R}}^{i+1}}^* q$ , *i.e.* how to normalize the transition  $r\sigma \rightarrow q$ ? An efficient matching algorithm based on tree automata intersection is described in [10, 12]. Normalization is described in the next section.

## 4.2 Normalization

The normalization function normalizes subterms by either states of  $\mathcal{Q}$  (using transitions of  $\Delta$ ) or new states. A state  $q$  of  $\mathcal{Q}$  is used to normalize a term  $t$  if  $t \rightarrow_{\Delta}^{\not\in} q$ . Normalizing by reusing states of  $\mathcal{Q}$  and transitions of  $\Delta$  permits to preserve the  $\not\in$ -determinism of  $\rightarrow_{\Delta}^{\not\in}$ . Indeed,  $\rightarrow_{\Delta}^{\not\in}$  can be kept deterministic during completion though  $\rightarrow_{\Delta}$  cannot.

**Definition 1 (New state).** *Given a set of transitions  $\Delta$ , a new state (for  $\Delta$ ) is a state of  $\mathcal{Q}$  not occurring in any left or right-hand side of any rule of  $\Delta$ <sup>1</sup>.*

We here define normalization as a bottom-up process. This definition is simpler and equivalent to top-down definitions [15]. In the recursive call, the choice of the context  $C[\ ]$  may be non deterministic but all the possible results are the equivalent modulo state renaming.

**Definition 2 (Normalization).** *Given a set of transitions  $\Delta$  defined on a set of states  $\mathcal{Q}$ , the normalization operation takes a transition  $t \rightarrow q$  such that  $t \in \mathcal{T}(\mathcal{F} \cup \mathcal{Q}) \setminus \mathcal{Q}$ , and  $q$  a new state for  $\Delta$ . Let  $C[\ ]$  be a non empty context of  $\mathcal{T}(\mathcal{F} \cup \mathcal{Q}) \setminus \mathcal{Q}$ ,  $f \in \mathcal{F}$  of arity  $n$  and  $q_1, \dots, q_n \in \mathcal{Q}$ . Normalization inductively generates a set of normalized transitions, by applying the following rules:*

1.  $Norm_{\Delta}(f(q_1, \dots, q_n) \rightarrow q) = \{f(q_1, \dots, q_n) \rightarrow q\}$
2.  $Norm_{\Delta}(C[f(q_1, \dots, q_n)] \rightarrow q) = \{f(q_1, \dots, q_n) \rightarrow q'\} \cup Norm_{\Delta \cup \{f(q_1, \dots, q_n) \rightarrow q'\}}(C[q'] \rightarrow q)$   
*where  $f(q_1, \dots, q_n) \rightarrow q' \in \Delta$  or,  $q'$  is a new state for  $\Delta$  and  $\forall q'' \in \mathcal{Q} : f(q_1, \dots, q_n) \rightarrow q'' \notin \Delta$ .*

We illustrate the above definition on the normalization of a simple transition.

*Example 1.* Given  $\Delta = \{b \rightarrow q_0\}$ ,  $Norm_{\Delta}(f(g(a), b, g(a)) \rightarrow q) = \{a \rightarrow q_1, g(q_1) \rightarrow q_2, b \rightarrow q_0, f(q_2, q_0, q_2) \rightarrow q\}$

**Lemma 1 ( $Norm_{\Delta}$  respects transitions of  $\Delta$ ).** *Let  $\Delta$  be an  $\not\in$ -deterministic set of transitions,  $q_{new}$  a new state for  $\Delta$  and  $t \in \mathcal{T}(\mathcal{F} \cup \mathcal{Q})$  s.t. there exists no state  $q'$  such that  $t \rightarrow_{\Delta}^{\not\in} q'$ . If  $f(q_1, \dots, q_n) \rightarrow q \in \Delta$  and  $f(q_1, \dots, q_n) \rightarrow q' \in Norm_{\Delta}(t \rightarrow q_{new})$  then  $q = q'$ .*

<sup>1</sup> Since  $\mathcal{Q}$  is a countably infinite set of states and  $\Delta$  is finite, a new state can always be found.

*Proof.* The first thing to remark is that the “ $\Delta$  parameter” of the *Norm* function only increases and remains  $\not\Leftarrow$ -deterministic, whatever the recursive calls may be, if its initial value is. This is a consequence of case 2 of the Definition 2 where we add to this parameter the transition  $f(q_1, \dots, q_n) \rightarrow q'$  only if there exists no transition  $f(q_1, \dots, q_n) \rightarrow q''$  in  $\Delta$ .

Now, we assume that  $f(q_1, \dots, q_n) \rightarrow q \in \Delta$ , and we prove that if there is a transition  $f(q_1, \dots, q_n) \rightarrow q' \in \text{Norm}_\Delta(t \rightarrow q_{\text{new}})$  then  $q = q'$ . The proof is done by induction on the number of symbols (of  $\mathcal{F}$ ) in  $t$ . If  $t$  has one symbol then it is of the form  $g(q'_1, \dots, q'_m)$ . We can only apply the case 1 of the definition. If  $t \neq f(q_1, \dots, q_n)$  then the result is  $\{t \rightarrow q\}$  and the property trivially holds. The other situation where  $t = f(q_1, \dots, q_n)$  is not possible since  $f(q_1, \dots, q_n) \rightarrow q \in \Delta$  contradicts the assumption  $t = f(q_1, \dots, q_n) \not\Leftarrow_\Delta^* q$ .

Now, assume that the property is true for  $t$  whose number of symbols is lesser or equal to  $n$ . For  $f(q_1, \dots, q_n) \rightarrow q'$  to belong to  $\text{Norm}_\Delta(t \rightarrow q_{\text{new}})$  it is necessarily added by case 2 of the definition of *Norm*. Hence, there exists a recursive call to *Norm* of the form  $\text{Norm}_{\Delta'}(C[f(q_1, \dots, q_n)] \rightarrow q_{\text{new}})$  where  $C[\ ]$  is a non empty context and  $\Delta' \supseteq \Delta$ . Since the transition  $f(q_1, \dots, q_n) \rightarrow q$  is in  $\Delta$  then it is in  $\Delta'$  and, by Definition 2, the added transition will be  $f(q_1, \dots, q_n) \rightarrow q$  (i.e.  $q = q'$ ). Thus, as explained above, we know that  $\Delta'$  is  $\not\Leftarrow$ -deterministic and that  $f(q_1, \dots, q_n) \rightarrow q \in \Delta'$ . Thus  $\Delta' \cup \{f(q_1, \dots, q_n) \rightarrow q\}$  is  $\not\Leftarrow$ -deterministic and we can use the induction hypothesis on  $\text{Norm}_{\Delta \cup \{f(q_1, \dots, q_n) \rightarrow q\}}(C[q] \rightarrow q_{\text{new}})$  and obtain that if  $f(q_1, \dots, q_n) \rightarrow q'' \in \text{Norm}_{\Delta \cup \{f(q_1, \dots, q_n) \rightarrow q\}}(C[q] \rightarrow q)$  then  $q = q''$ .

**Lemma 2 (Result of  $\text{Norm}_\Delta$  is  $\not\Leftarrow$ -deterministic).** *Let  $\Delta$  be an  $\not\Leftarrow$ -deterministic set of transitions,  $q_{\text{new}}$  a new state for  $\Delta$  and  $t \in \mathcal{T}(\mathcal{F} \cup \mathcal{Q})$  s.t. there exists no state  $q'$  such that  $t \rightarrow_\Delta^* q'$ . The set  $\text{Norm}_\Delta(t \rightarrow q_{\text{new}})$  is  $\not\Leftarrow$ -deterministic.*

*Proof.* We show this lemma by contradiction. Assume that  $\text{Norm}_\Delta(t \rightarrow q_{\text{new}})$  is not  $\not\Leftarrow$ -deterministic. Thus, there exists a configuration  $c \in \mathcal{T}(\mathcal{F} \cup \mathcal{Q}) \setminus \mathcal{Q}$  and two states  $q, q'$  such that  $q \neq q'$  and  $\{c \rightarrow q, c \rightarrow q'\} \subseteq \text{Norm}_\Delta(t \rightarrow q_{\text{new}})$ . Since there are (at least) two transitions in  $\text{Norm}_\Delta(t \rightarrow q_{\text{new}})$ , we know that (at least) one transition in  $\{c \rightarrow q, c' \rightarrow q\}$  have been added by the case 2 of Definition 2. Now, let  $c = f(q_1, \dots, q_n)$ . Assume that  $c \rightarrow q$  is found in recursive calls of *Norm* before  $c \rightarrow q'$ . The recursive call is thus of the form:  $\{f(q_1, \dots, q_n) \rightarrow q\} \cup \text{Norm}_{\Delta' \cup \{f(q_1, \dots, q_n) \rightarrow q\}}(C[q] \rightarrow q_{\text{new}})$ , where  $\Delta' \supseteq \Delta$ . Furthermore since  $c \rightarrow q'$  is not in  $\Delta'$ , it is in  $\text{Norm}_{\Delta' \cup \{f(q_1, \dots, q_n) \rightarrow q\}}(C[q] \rightarrow q_{\text{new}})$ . However, using Lemma 1, we obtain that  $q = q'$  which is a contradiction.

### 4.3 One step of completion

A step of completion only consists in joining critical pairs. We first need to formally define the substitutions under concern: *regular language substitutions*.

**Definition 3 (Regular language substitution).** *A regular language substitution over an automaton  $\mathcal{A}$  with a set of states  $\mathcal{Q}$  is a function  $\sigma : \mathcal{X} \mapsto \mathcal{Q}$ . We can extend this definition to a morphism  $\sigma : \mathcal{T}(\mathcal{F}, \mathcal{X}) \mapsto \mathcal{T}(\mathcal{F}, \mathcal{Q})$ . We denote by  $\Sigma(\mathcal{Q}, \mathcal{X})$  the set of regular language substitutions built over  $\mathcal{Q}$  and  $\mathcal{X}$ .*

**Definition 4 (Set of critical pairs).** *Let a TRS  $\mathcal{R}$  and a tree automaton  $\mathcal{A} = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_f, \Delta \rangle$ . The set of critical pairs between  $\mathcal{R}$  and  $\mathcal{A}$  is  $CP(\mathcal{R}, \mathcal{A}) = \{(l \rightarrow r, \sigma, q) \mid l \rightarrow r \in \mathcal{R}, q \in \mathcal{Q}, \sigma \in \Sigma(\mathcal{Q}, \mathcal{X}), l\sigma \rightarrow_{\mathcal{A}}^* q, r\sigma \not\Leftarrow_{\mathcal{A}}^* q \text{ and } (l \rightarrow r, \sigma, q)\}$ .*

Recall that the completion process will build a sequence  $\mathcal{A}_{\mathcal{R}}^0, \mathcal{A}_{\mathcal{R}}^1, \dots, \mathcal{A}_{\mathcal{R}}^k$  of automata such that if  $s \in \mathcal{L}(\mathcal{A}_{\mathcal{R}}^i)$  and  $s \rightarrow_{\mathcal{R}} t$  then  $t \in \mathcal{L}(\mathcal{A}_{\mathcal{R}}^{i+1})$ . One step of completion, i.e. the process computing  $\mathcal{A}_{\mathcal{R}}^{i+1}$  from  $\mathcal{A}_{\mathcal{R}}^i$ , is defined as follows. Again, the following definition is a simplification of the definition of [15].

**Definition 5 (One step automaton completion).** *Let  $\mathcal{A} = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_f, \Delta \rangle$  be a tree automaton,  $\mathcal{R}$  be a left-linear TRS. The one step completed automaton is  $\mathcal{C}_{\mathcal{R}}(\mathcal{A}) = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_f, \text{Join}^{CP(\mathcal{R}, \mathcal{A})}(\Delta) \rangle$  where  $\text{Join}^S(\Delta)$  is inductively defined by:*

- $Join^\emptyset(\Delta) = \Delta$
- $Join^{\{(l \rightarrow r, q, \sigma)\} \cup S}(\Delta) = Join^S(\Delta \cup \Delta')$  where
  - $\Delta' = \{q' \rightarrow q\}$  if there exists  $q' \in \mathcal{Q}$  s.t.  $r\sigma \rightarrow_{\Delta}^{q'} q'$ , and otherwise
  - $\Delta' = Norm_\Delta(r\sigma \rightarrow q') \cup \{q' \rightarrow q\}$  where  $q'$  is a new state for  $\Delta$

*Example 2.* Let  $\mathcal{A}$  be a tree automaton with  $\Delta = \{f(q_1) \rightarrow q_0, a \rightarrow q_1, g(q_1) \rightarrow q_2\}$ .

- If  $\mathcal{R} = \{f(a) \rightarrow g(a)\}$  then  $CP(\mathcal{R}, \mathcal{A}) = \{(f(a) \rightarrow g(a), \sigma_1, q_0)\}$  with  $\sigma_1 = \emptyset$  because  $f(a)\sigma_1 \rightarrow_{\mathcal{A}}^* q_0$  and  $f(a)\sigma_1 \rightarrow_{\mathcal{R}} g(a)\sigma_1$ . Besides, we have  $g(a) \rightarrow_{\mathcal{A}}^{q'} q_2$ . Hence  $Join^{\{(f(a) \rightarrow g(a), \sigma_1, q_0)\}}(\Delta) = Join^\emptyset(\Delta \cup \{q_2 \rightarrow q_0\}) = \Delta \cup \{q_2 \rightarrow q_0\}$ . Thus,  $\mathcal{C}_{\mathcal{R}}(\mathcal{A}) = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_f, \Delta \cup \{q_2 \rightarrow q_0\} \rangle$ ;
- If  $\mathcal{R} = \{f(x) \rightarrow x\}$  then  $CP(\mathcal{R}, \mathcal{A}) = \{(f(x) \rightarrow x, \sigma_2, q_0)\}$  with  $\sigma_2 = \{x \mapsto q_1\}$  since  $f(x)\sigma_2 \rightarrow_{\mathcal{A}}^* q_0$  and  $f(x)\sigma_2 \rightarrow_{\mathcal{R}} x\sigma_2 = q_1$ . Similarly, from  $q_1 \rightarrow_{\mathcal{A}}^{q'} q_1$  we get that  $Join^{\{(f(x) \rightarrow x, \sigma_2, q_0)\}}(\Delta) = Join^\emptyset(\Delta \cup \{q_1 \rightarrow q_0\}) = \Delta \cup \{q_1 \rightarrow q_0\}$  and  $\mathcal{C}_{\mathcal{R}}(\mathcal{A}) = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_f, \Delta \cup \{q_1 \rightarrow q_0\} \rangle$ ;
- If  $\mathcal{R} = \{f(x) \rightarrow f(g(x))\}$  then  $CP(\mathcal{R}, \mathcal{A}) = \{(f(x) \rightarrow f(g(x)), \sigma_3, q_0)\}$  with  $\sigma_3 = \{x \mapsto q_1\}$ , because  $f(x)\sigma_3 \rightarrow_{\mathcal{A}}^* q_0$  and  $f(x)\sigma_3 \rightarrow_{\mathcal{R}} f(g(x))\sigma_3$ . We have  $f(g(x))\sigma_3 = f(g(q_1))$  and this time, there exists no state  $q$  such that  $f(g(q_1)) \rightarrow_{\mathcal{A}}^{q'} q$ . Hence,  $Join^{\{(f(x) \rightarrow f(g(x)), \sigma_3, q_0)\}}(\Delta) = Join^\emptyset(\Delta \cup Norm_\Delta(f(g(q_1)) \rightarrow q_3) \cup \{q_3 \rightarrow q_0\})$ . Since  $Norm_\Delta(f(g(q_1)) \rightarrow q_3) = \{f(q_2) \rightarrow q_3, q(q_1) \rightarrow q_2\}$ , we get that  $\mathcal{C}_{\mathcal{R}}(\mathcal{A}) = \langle \mathcal{F}, \mathcal{Q} \cup \{q_3\}, \mathcal{Q}_f, \Delta \cup \{f(q_2) \rightarrow q_3, q_3 \rightarrow q_0\} \rangle$ .

#### 4.4 Simplification of Tree Automata by Equations

In this section, we define the *simplification* of tree automata  $\mathcal{A}$  w.r.t. a set of equations  $E$ . This operation will be necessary to over-approximate languages that cannot be recognized *exactly* using tree automata completion, *e.g.* non regular languages. The simplification operation consists in finding  $E$ -equivalent terms recognized in  $\mathcal{A}$  by different states and then by merging those states together. The merging of states is performed using renaming of a state in a tree automaton.

**Definition 6 (Renaming of a state in a tree automaton).** Let  $\mathcal{Q}, \mathcal{Q}'$  be set of states,  $\mathcal{A} = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_f, \Delta \rangle$  be a tree automaton, and  $\alpha$  a function  $\alpha : \mathcal{Q} \mapsto \mathcal{Q}'$ . We denote by  $\mathcal{A}\alpha$  the tree automaton where every occurrence of  $q$  is replaced by  $\alpha(q)$  in  $\mathcal{Q}$ ,  $\mathcal{Q}_f$  and in every left and right-hand side of every transition of  $\Delta$ .

If there exists a bijection  $\alpha$  such that  $\mathcal{A} = \mathcal{A}'\alpha$  then  $\mathcal{A}$  and  $\mathcal{A}'$  are said to be *equivalent modulo renaming*. Now we define the *simplification relation* which merges states in a tree automaton according to an equation. Note that it is not required for equations of  $E$  to be linear.

**Definition 7 (Simplification relation).** Let  $\mathcal{A} = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_f, \Delta \rangle$  be a tree automaton and  $E$  be a set of equations. For  $s = t \in E$ ,  $\sigma \in \Sigma(\mathcal{Q}, \mathcal{X})$ ,  $q_a, q_b \in \mathcal{Q}$  such that  $s\sigma \rightarrow_{\mathcal{A}}^{q'} q_a$ ,  $t\sigma \rightarrow_{\mathcal{A}}^{q'} q_b$ , *i.e.*

$$\begin{array}{ccc} s\sigma & \xlongequal{E} & t\sigma \\ \mathcal{A}, \sigma' \downarrow * & & * \downarrow \mathcal{A}, \sigma' \\ q_a & & q_b \end{array}$$

and  $q_a \neq q_b$  then  $\mathcal{A}$  can be simplified into  $\mathcal{A}' = \mathcal{A}\{q_b \mapsto q_a\}$ , denoted by  $\mathcal{A} \rightsquigarrow_E \mathcal{A}'$ .  $\diamond$

*Example 3.* Let  $E = \{s(s(x)) = s(x), a = b\}$  and  $\mathcal{A}$  be the tree automaton with  $\mathcal{Q}_f = \{q_2, q_4\}$  and set of transitions  $\Delta = \{a \rightarrow q_0, s(q_0) \rightarrow q_1, s(q_1) \rightarrow q_2, b \rightarrow q_3, s(q_3) \rightarrow q_4\}$ . Hence  $\mathcal{L}(\mathcal{A}) = \{s(s(a)), s(b)\}$ . We can perform a first simplification step using the equation  $s(s(x)) = s(x)^2$ ,

<sup>2</sup> Note that we could have begun to simplify  $\mathcal{A}$  w.r.t. equation  $a = b$ , but as we will see below, this makes no difference.

because we found a substitution  $\sigma = \{x \mapsto q_0\}$  such that:

$$\begin{array}{ccc} s(s(q_0)) & \xlongequal{E} & s(q_0) \\ \mathcal{A}, \epsilon' \downarrow * & & * \downarrow \mathcal{A}, \epsilon' \\ q_2 & & q_1 \end{array}$$

Hence,  $\mathcal{A} \sim_E \mathcal{A}' = \mathcal{A}\{q_2 \mapsto q_1\}$ <sup>3</sup> Thus,  $\mathcal{A}'$  is the automaton with  $\mathcal{Q}'_f = \{q_1, q_4\}$  and  $\Delta = \{a \rightarrow q_0, s(q_0) \rightarrow q_1, s(q_1) \rightarrow q_1, b \rightarrow q_3, s(q_3) \rightarrow q_4\}$  and  $\mathcal{L}(\mathcal{A}') = \{s^*(s(a)), s(b)\}$ . Then, we can perform a second simplification step using the equation  $a = b$ , because we found a substitution  $\sigma' = \emptyset$  such that:

$$\begin{array}{ccc} a & \xlongequal{E} & b \\ \mathcal{A}, \epsilon' \downarrow * & & * \downarrow \mathcal{A}, \epsilon' \\ q_0 & & q_3 \end{array}$$

We can thus simplify  $\mathcal{A}'$  in this way:  $\mathcal{A}' \sim_E \mathcal{A}'' = \mathcal{A}'\{q_0 \mapsto q_3\}$  where  $\mathcal{A}''$  is the tree automaton such that  $\mathcal{Q}''_f = \mathcal{Q}'_f$  and  $\Delta'' = \{a \rightarrow q_3, s(q_3) \rightarrow q_1, s(q_1) \rightarrow q_1, b \rightarrow q_3, s(q_3) \rightarrow q_4\}$ . A last step of simplification can be performed using  $s(s(x)) = s(x)$  and leads to the automaton  $\mathcal{A}''' = \mathcal{A}''\{q_4 \mapsto q_1\}$  with  $\mathcal{Q}'''_f = \{q_1\}$  and  $\Delta''' = \{a \rightarrow q_3, s(q_3) \rightarrow q_1, s(q_1) \rightarrow q_1, b \rightarrow q_3\}$ . Automaton  $\mathcal{A}'''$  cannot be simplified, is thus a normal form of  $\sim_E$  and  $\mathcal{L}(\mathcal{A}''') = \{s^*(s(a|b))\}$ .

As stated in [15] and to no one's surprise, simplification  $\sim_E$  is a terminating relation (each step suppresses a state) and it enlarges the language recognized by a tree automaton, i.e. if  $\mathcal{A} \sim_E \mathcal{A}'$  then  $\mathcal{L}(\mathcal{A}) \subseteq \mathcal{L}(\mathcal{A}')$ . More surprisingly, no matter how simplification steps are performed, there exists a canonical simplified tree automaton. In the following,  $\mathcal{A} \sim_E^! \mathcal{A}'$  denotes that  $\mathcal{A} \sim_E^* \mathcal{A}'$  and  $\mathcal{A}'$  is irreducible by  $\sim_E$ , i.e. no simplification by  $E$  can be performed on  $\mathcal{A}'$ .

**Theorem 1 (Canonical Simplified Tree Automata [15]).** *Let  $\mathcal{A}, \mathcal{A}'_1, \mathcal{A}'_2$  be tree automata and  $E$  be a set of equations. If  $\mathcal{A} \sim_E^! \mathcal{A}'_1$  and  $\mathcal{A} \sim_E^! \mathcal{A}'_2$  then  $\mathcal{A}'_1$  and  $\mathcal{A}'_2$  are equivalent modulo a bijective renaming.*

In the following, we note  $\mathcal{S}_E(\mathcal{A})$  this canonical simplified automaton, i.e. one of the possible automaton  $\mathcal{A}'$  such that  $\mathcal{A} \sim_E^! \mathcal{A}'$ .

#### 4.5 The full Completion Algorithm

Now, we can define the full equational completion algorithm.

**Definition 8 (Automaton completion).** *Let  $\mathcal{A}$  be a tree automaton,  $\mathcal{R}$  a left-linear TRS and  $E$  a set of equations.*

- $\mathcal{A}_{\mathcal{R}, E}^0 = \mathcal{A}$ ,
- $\mathcal{A}_{\mathcal{R}, E}^{n+1} = \mathcal{S}_E(\mathcal{C}_{\mathcal{R}}(\mathcal{A}_{\mathcal{R}, E}^n))$
- $\mathcal{A}_{\mathcal{R}, E}^*$  is a fixpoint, i.e.  $\mathcal{A}_{\mathcal{R}, E}^* = \mathcal{A}_{\mathcal{R}, E}^k = \mathcal{A}_{\mathcal{R}, E}^{k+1}$  with  $k \in \mathbb{N}$ .

In practice, a good criterion to know that  $\mathcal{A}_{\mathcal{R}, E}^k$  is a fixpoint is when  $CP(\mathcal{R}, \mathcal{A}_{\mathcal{R}, E}^k) = \emptyset$ . However, a fixpoint cannot always be finitely reached<sup>4</sup>. Another way to ensure termination is to provide a set of approximating equations that is able to overcome infinite rewriting and thus completion divergence.

<sup>3</sup> or  $\{q_1 \mapsto q_2\}$ , any of  $q_1$  or  $q_2$  can be used for renaming.

<sup>4</sup> See [12], for classes of  $\mathcal{R}$  for which a fixpoint always exists.



*Example 4.* Let  $\mathcal{R} = \{f(x, y) \rightarrow f(s(x), s(y))\}$ ,  $E = \{s(s(x)) = s(x)\}$  and  $\mathcal{A}^0$  be the tree automaton with set of transitions  $\Delta = \{f(q_a, q_b) \rightarrow q_0, a \rightarrow q_a, b \rightarrow q_b\}$ , i.e.  $\mathcal{L}(\mathcal{A}^0) = \{f(a, b)\}$ . The completion ends after two completion steps on  $\mathcal{A}_{\mathcal{R}, E}^2$  which is a fixpoint. Completion steps are summed up in the following table. To simplify the presentation, we do not repeat the common transitions:  $\mathcal{A}_{\mathcal{R}, E}^i$  and  $\mathcal{C}_{\mathcal{R}}(\mathcal{A}^i)$  columns are supposed to contain all transitions of  $\mathcal{A}^0, \dots, \mathcal{A}_{\mathcal{R}, E}^{i-1}$ .

$\mathcal{A}^0$	$\mathcal{C}_{\mathcal{R}}(\mathcal{A}^0)$	$\mathcal{A}_{\mathcal{R}, E}^1$	$\mathcal{C}_{\mathcal{R}}(\mathcal{A}_{\mathcal{R}, E}^1)$	$\mathcal{A}_{\mathcal{R}, E}^2$
$f(q_a, q_b) \rightarrow q_0$	$f(q_1, q_2) \rightarrow q_3$	$f(q_1, q_2) \rightarrow q_3$	$f(q_4, q_5) \rightarrow q_6$	$f(q_1, q_2) \rightarrow q_6$
$a \rightarrow q_a$	$s(q_a) \rightarrow q_1$	$s(q_a) \rightarrow q_1$	$s(q_1) \rightarrow q_4$	$s(q_1) \rightarrow q_1$
$b \rightarrow q_b$	$s(q_b) \rightarrow q_2$	$s(q_b) \rightarrow q_2$	$s(q_2) \rightarrow q_5$	$s(q_2) \rightarrow q_2$
	$q_3 \rightarrow q_0$	$q_3 \rightarrow q_0$	$q_6 \rightarrow q_3$	

The automaton  $\mathcal{A}_{\mathcal{R}, E}^1$  is exactly  $\mathcal{C}_{\mathcal{R}}(\mathcal{A}^0)$  since simplification by equations do not apply. Then  $\mathcal{C}_{\mathcal{R}}(\mathcal{A}_{\mathcal{R}, E}^1)$  contains all the transitions of  $\mathcal{A}_{\mathcal{R}, E}^1$  plus those obtained by the resolution of the critical pair  $f(q_1, q_2) \rightarrow_{\mathcal{A}^*} q_3$  and  $f(q_1, q_2) \rightarrow_{\mathcal{R}} f(s(q_1), s(q_2))$ . On  $\mathcal{C}_{\mathcal{R}}(\mathcal{A}_{\mathcal{R}, E}^1)$  simplification using the equation  $s(s(x)) = s(x)$  can be applied on following instances:  $s(s(q_a)) = s(q_a)$  and  $s(s(q_b)) = q_b$  which results in merging states  $q_4$  with  $q_1$  and  $q_5$  with  $q_2$ . Thus,  $\mathcal{A}_{\mathcal{R}, E}^2 = \mathcal{C}_{\mathcal{R}}(\mathcal{A}_{\mathcal{R}, E}^1) \setminus \{q_4 \mapsto q_1, q_5 \mapsto q_2\}$ . This last automaton is a fixed point because  $CP(\mathcal{R}, \mathcal{A}_{\mathcal{R}, E}^2) = \emptyset$ .

Now, we recall the lower and upper bound theorems. Tree automata completion of automaton  $\mathcal{A}$  with TRS  $\mathcal{R}$  and set of equations  $E$  is lower bounded by  $\mathcal{R}^*(\mathcal{L}(\mathcal{A}))$  and upper bounded by  $\mathcal{R}_E^*(\mathcal{L}(\mathcal{A}))$ . The lower bound theorem ensures that the completed automaton  $\mathcal{A}_{\mathcal{R}, E}^*$  recognizes all  $\mathcal{R}$ -reachable terms (but not all  $\mathcal{R}/E$ -reachable terms). The upper bound theorem guarantees that all terms recognized by  $\mathcal{A}_{\mathcal{R}, E}^*$  are only  $\mathcal{R}/E$ -reachable terms.

**Theorem 2 (Lower bound [15]).** *Let  $\mathcal{R}$  be a left-linear TRS,  $\mathcal{A}$  be a tree automaton and  $E$  be a set of equations. If completion terminates on  $\mathcal{A}_{\mathcal{R}, E}^*$  then*

$$\mathcal{L}(\mathcal{A}_{\mathcal{R}, E}^*) \supseteq \mathcal{R}^*(\mathcal{L}(\mathcal{A}))$$

Note that the left-linearity condition on  $\mathcal{R}$  can be removed using, so-called, packed states [27]. This condition can also be weakened using the left-linearity condition [10] or conditions on languages matched by non linear variables [4]. The upper bound theorem states the precision result of completion. It is defined using the  $\mathcal{R}/E$ -coherence property. The intuition behind  $\mathcal{R}/E$ -coherence is the following: in the tree automaton  $\epsilon$ -transitions represent rewriting steps and normalized transitions recognize  $E$ -equivalence classes. More precisely, in a  $\mathcal{R}/E$ -coherent tree automaton, if two terms  $s, t$  are recognized into the same state  $q$  using only normalized transitions then they belong to the same  $E$ -equivalence class. Otherwise, if at least one  $\epsilon$ -transition is necessary to recognize, say,  $t$  into  $q$  then at least one step of rewriting was necessary to obtain  $t$  from  $s$ .

**Theorem 3 (Upper bound [15]).** *Let  $\mathcal{R}$  be a left-linear TRS,  $E$  a set of equations and  $\mathcal{A}$  a  $\mathcal{R}/E$ -coherent tree automaton. For any  $i \in \mathbb{N}$ :*

$$\mathcal{L}(\mathcal{A}_{\mathcal{R}, E}^i) \subseteq \mathcal{R}_E^*(\mathcal{L}(\mathcal{A})) \quad \text{and} \quad \mathcal{A}_{\mathcal{R}, E}^i \text{ is } \mathcal{R}/E\text{-coherent}$$

The fact that those two theorems apply on different sets, namely  $\mathcal{R}^*$  and  $\mathcal{R}_E^*$  is important to use this technique for software verification. Indeed, if  $\mathcal{R}$  models the program and  $E$  defines the approximation then it is natural to focus the theorem on the over-approximation of  $\mathcal{R}$ -reachable terms rather than on  $\mathcal{R}/E$ -reachable ones. In the context of verification,  $\mathcal{R}/E$ -reachable terms that are not  $\mathcal{R}$ -reachable are not interesting: they are necessarily part of the approximation defined by  $E$ . Computing exactly or over-approximating  $\mathcal{R}/E$ -reachable terms is nevertheless possible for some well identified classes of  $E$  [15].

## 5 Termination criterion for a given set of equations

Given a set of equations  $E$ , the effect of the simplification with  $E$  on a tree automaton is to merge two distinct states recognizing instances of the left and right-hand side for all the equations of  $E$ . In this section, we give a sufficient condition on  $E$  and on the completed tree automata  $\mathcal{A}_{\mathcal{R},E}^1$  for the tree automata completion to always terminate. The intuition behind this condition is simple: if the set of equivalence classes for  $E$ , *i.e.*  $\mathcal{T}(\mathcal{F})/\equiv_E$ , is finite then so should be the set of new states used in completion. However, this is not true in general because simplification of an automaton with  $E$  does not necessarily merge all  $E$ -equivalent terms.

*Example 5.* Let  $\mathcal{A}$  be the tree automaton with set of transitions  $a \rightarrow q$ ,  $\mathcal{R} = \{a \rightarrow c\}$  and let  $E = \{a = b, b = c\}$ . The set of transitions of  $\mathcal{C}_{\mathcal{R}}(\mathcal{A})$  is  $\{a \rightarrow q, c \rightarrow q', q' \rightarrow q\}$ . We have  $a \equiv_E c$ ,  $a \in \mathcal{L}^{\mathcal{C}}(\mathcal{C}_{\mathcal{R}}(\mathcal{A}), q)$  and  $c \in \mathcal{L}^{\mathcal{C}}(\mathcal{C}_{\mathcal{R}}(\mathcal{A}), q')$  but on the automaton  $\mathcal{C}_{\mathcal{R}}(\mathcal{A})$ , no simplification situation (as described by Definition 7), can be found because the term  $b$  is not recognized by  $\mathcal{C}_{\mathcal{R}}(\mathcal{A})$ . Hence, the simplified automaton is  $\mathcal{C}_{\mathcal{R}}(\mathcal{A})$  where  $a$  and  $c$  are recognized by different states.

There is no simple solution to have a simplification algorithm merging all states recognizing  $E$ -equivalent terms (see Section 6). Besides, having an automaton  $\mathcal{A}$  that is complete before completion would apparently solve the problem. This is the case for the above example. If for all term  $t \in \mathcal{T}(\mathcal{F})$  there exists at least a state  $q$  such that  $t \in \mathcal{L}^{\mathcal{C}}(\mathcal{A}, q)$  then we would have a transition  $b \rightarrow q''$  and thus simplification could have been performed until having  $a$  and  $c$  recognized by the same state. However, using *complete* initial automata to compute over-approximation of reachable may produce very rough approximations. This is the case when the structure of the complete initial tree automaton interfere with  $E$  and completion thus add transitions recognizing unreachable terms in final states.

*Example 6.* Let  $\mathcal{F} = \{a, b, c\}$ ,  $\mathcal{R} = \{a \rightarrow b\}$ ,  $E = \{b = c\}$  and  $\mathcal{A}$  the complete tree automaton with  $\mathcal{Q}_{\mathcal{F}} = \{q_0\}$  and  $\Delta = \{a \rightarrow q_0, b \rightarrow q_1, c \rightarrow q_1\}$ . The first completion step yields the transition  $q_1 \rightarrow q_0$ . The transition set of the final automaton  $\mathcal{A}_{\mathcal{R},E}^1$  is thus  $\{a \rightarrow q_0, q_1 \rightarrow q_0, b \rightarrow q_1, c \rightarrow q_1\}$   $\mathcal{L}(\mathcal{A}_{\mathcal{R},E}^1) = \{a, b, c\}$  which is a coarse approximation of  $\mathcal{R}^*(\mathcal{L}(\mathcal{A}))$ . This result has to be compared with the result obtained when completing an equivalent initial tree automaton  $\mathcal{B}$  which is not complete. Let  $\Delta' = \{a \rightarrow q_0\}$  be the set of transitions of  $\mathcal{B}$ . Completion of  $\mathcal{B}$  stops on  $\mathcal{B}_{\mathcal{R},E}^1$  with transitions  $\{a \rightarrow q_0, b \rightarrow q'_1, q'_1 \rightarrow q_0\}$  where  $q'_1$  is a new state and  $\mathcal{L}(\mathcal{B}_{\mathcal{R},E}^1) = \{a, b\}$  which is precisely  $\mathcal{R}^*(\mathcal{L}(\mathcal{A}'))$  and equal to  $\mathcal{R}^*(\mathcal{L}(\mathcal{A}))$ .

In the next section, we propose to give some simple restrictions on  $E$  to ensure that completion terminates. In Section 5.2, we will see how those restrictions can easily be met for “functional” TRS, *i.e.* a typed first-order functional program translated into a TRS.

### 5.1 General criterion

What Example 5 shows is that, for a simplification with  $E$  to apply, it is necessary that both sides of the equation are recognized by the tree automaton. In the following, we will define a set  $E^c$  of *contracting* equations so that this property is true. What Example 5 does not show is that, by default, tree automata are not  $E$ -compatible. In particular, any non  $\epsilon$ -deterministic automaton does not satisfy the reflexivity of  $\equiv_E$ . For instance, if an automaton  $\mathcal{A}$  has two transitions  $a \rightarrow q_1$  and  $a \rightarrow q_2$ , since  $a \equiv_E a$  for all  $E$ , for  $\mathcal{A}$  to be  $E$ -compatible we should have  $q_1 = q_2$ . To enforce  $\epsilon$ -determinism by automata simplification, we define a set of *reflexivity equations* as follows.

**Definition 9 (Set of reflexivity equations  $E^r$ ).** For a given set of symbols  $\mathcal{F}$ ,  $E^r = \{f(x_1, \dots, x_n) = f(x_1, \dots, x_n) \mid f \in \mathcal{F} \text{ and arity of } f \text{ is } n\}$ .

Note that for all set of equations  $E$ , the relation  $\equiv_E$  is trivially equivalent to  $\equiv_{E \cup E^r}$ . Furthermore, simplification with  $E^r$  transforms all automaton into an  $\epsilon$ -deterministic automaton, as stated in the following lemma.

**Lemma 3.** *For all tree automaton  $\mathcal{A}$  and all set of equation  $E$ , if  $E \supseteq E^r$  and  $\mathcal{A} \sim_E^! \mathcal{A}'$  then  $\mathcal{A}'$  is  $\not\sim$ -deterministic.*

*Proof.* We prove this by induction on the height of the terms recognized by  $\mathcal{A}'$ . This is true for constants because otherwise there would be a constant  $a$  such that  $a \rightarrow_{\mathcal{A}'}^{\not\sim^*} q$  and  $a \rightarrow_{\mathcal{A}'}^{\not\sim^*} q'$  with  $q \neq q'$ . However since  $a = a \in E^r$  we can simplify  $\mathcal{A}'$  which contradicts the fact that  $\mathcal{A}'$  is in normal form w.r.t.  $\sim_E$ . For the inductive case, assume that there exists a term  $t = f(t_1, \dots, t_n)$  such that  $t \rightarrow_{\mathcal{A}'}^{\not\sim^*} q$  and  $t \rightarrow_{\mathcal{A}'}^{\not\sim^*} q'$  with  $q \neq q'$ . Using the induction hypothesis, we know that for each  $t_i$  for  $i = 1 \dots n$  there exists a unique state  $q_i$  such that  $t_i \rightarrow_{\mathcal{A}'}^{\not\sim^*} q_i$ . Hence,  $f(t_1, \dots, t_n) \rightarrow_{\mathcal{A}'}^{\not\sim^*} f(q_1, \dots, q_n)$  but  $f(q_1, \dots, q_n) \rightarrow_{\mathcal{A}'}^{\not\sim^*} q$  and  $f(q_1, \dots, q_n) \rightarrow_{\mathcal{A}'}^{\not\sim^*} q'$ . However, this is a simplification situation for the equation  $f(x_1, \dots, x_n) = f(x_1, \dots, x_n) \in E^r$  which contradicts the fact that  $\mathcal{A}'$  is in normal form for  $\sim_E$ .

We now define the set  $E_{\mathcal{K}}^c$  of contracting equations. This set is defined for a set of symbols  $\mathcal{K}$  which can be a subset of  $\mathcal{F}$ . This will be used later to restrict contracting equations to the subset of constructor symbols of  $\mathcal{F}$ .

**Definition 10 (Set  $E_{\mathcal{K}}^c$  of contracting equations for  $\mathcal{K}$ ).** *Let  $\mathcal{K} \subseteq \mathcal{F}$ . The set  $E_{\mathcal{K}}^c$  is a set of contracting equations for  $\mathcal{K}$  if all equations of  $E_{\mathcal{K}}^c$  are of the form  $u = u|_p$  with  $u \in \mathcal{T}(\mathcal{K}, \mathcal{X})$  a linear term,  $p \neq \Lambda$ , and if the set of normal forms of  $\mathcal{T}(\mathcal{K})$  w.r.t. the TRS  $\overrightarrow{E_{\mathcal{K}}^c} = \{u \rightarrow u|_p \mid u = u|_p \in E_{\mathcal{K}}^c\}$  is finite.*

Contracting equations, if defined on  $\mathcal{F}$ , define an upper bound on the number of states of a simplified automaton.

**Lemma 4.** *Let  $\mathcal{A}$  be a tree automaton,  $E$  be a set of equations such that  $E \supseteq E_{\mathcal{F}}^c \cup E^r$ . The simplified automaton  $\mathcal{S}_E(\mathcal{A})$  is an  $\not\sim$ -deterministic automaton having no more states than terms in  $\text{IRR}(\overrightarrow{E_{\mathcal{F}}^c})$ .*

*Proof.* First, assume for all state  $q$  of  $\mathcal{S}_E(\mathcal{A})$ ,  $\mathcal{L}^{\not\sim}(\mathcal{S}_E(\mathcal{A}), q) \cap \text{IRR}(\overrightarrow{E_{\mathcal{F}}^c}) = \emptyset$ . Then, for all term  $s$  such that  $s \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} q$ , we know that  $s$  is not in normal form w.r.t.  $\overrightarrow{E_{\mathcal{F}}^c}$ . As a result, the left-hand side of an equation of  $E_{\mathcal{F}}^c$  can be applied to  $s$ . This means that there exists an equation  $u = u|_p$ , a ground context  $C$  and a substitution  $\theta$  such that  $s = C[u\theta]$ . Furthermore, since  $s \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} q$ , we know that  $C[u\theta] \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} q$  and that there exists a state  $q'$  such that  $C[q'] \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} q$  and  $u\theta \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} q'$ . From  $u\theta \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} q'$ , we know that all subterms of  $u\theta$  are recognized by at least one state in  $\mathcal{S}_E(\mathcal{A})$ . Thus, there exists a state  $q''$  such that  $u|_p\theta \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} q''$ . We thus have a situation of application of the equation  $u = u|_p$  in the automaton. Since  $\mathcal{S}_E(\mathcal{A})$  is simplified, we thus know that  $q' = q''$ . As mentioned above, we know that  $C[q'] \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} q$ . Hence  $C[u|_p\theta] \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} C[q'] \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} q$ . If  $C[u|_p\theta]$  is not in normal form w.r.t.  $\overrightarrow{E_{\mathcal{F}}^c}$  then we can do the same reasoning on  $C[u|_p\theta] \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} q$  until getting a term that is in normal form w.r.t.  $\overrightarrow{E_{\mathcal{F}}^c}$  and recognized by the same state  $q$ . Thus, this contradicts the fact that  $\mathcal{S}_E(\mathcal{A})$  recognizes no term of  $\text{IRR}(\overrightarrow{E_{\mathcal{F}}^c})$ .

Then, by definition of  $E_{\mathcal{F}}^c$ ,  $\text{IRR}(\overrightarrow{E_{\mathcal{F}}^c})$  is finite. Let  $\{t_1, \dots, t_n\}$  be the subset of  $\text{IRR}(\overrightarrow{E_{\mathcal{F}}^c})$  recognized by  $\mathcal{S}_E(\mathcal{A})$ . Let  $q_1, \dots, q_n$  be the states recognizing  $t_1, \dots, t_n$  respectively. We know that there is a finite set of states recognizing  $t_1, \dots, t_n$  because  $E \supseteq E^r$  and Lemma 3 entails that  $\mathcal{S}_E(\mathcal{A})$  is  $\not\sim$ -deterministic. Now, for all term  $s$  recognized by a state  $q$  in  $\mathcal{S}_E(\mathcal{A})$ , i.e.  $s \rightarrow_{\mathcal{S}_E(\mathcal{A})}^{\not\sim^*} q$ , we can use a reasoning similar to the one carried out above and show that  $q$  is equal to one state of  $\{q_1, \dots, q_n\}$  recognizing normal forms of  $\overrightarrow{E_{\mathcal{F}}^c}$  in  $\mathcal{S}_E(\mathcal{A})$ . Finally, there are at most  $\text{card}(\text{IRR}(\overrightarrow{E_{\mathcal{F}}^c}))$  states in  $\mathcal{S}_E(\mathcal{A})$ .

Now it is possible to state the Theorem guaranteeing the termination of completion if the set of equations  $E$  contains a set of contracting equations  $E_{\mathcal{F}}^c$  for  $\mathcal{F}$  and a set of reflexivity equations.

**Theorem 4.** *Let  $\mathcal{A}$  be a tree automaton,  $\mathcal{R}$  a left linear TRS and  $E$  a set of equations. If  $E \supseteq E^r \cup E_{\mathcal{F}}^c$  then completion of  $\mathcal{A}$  by  $\mathcal{R}$  and  $E$  terminates.*

*Proof.* For completion to diverge it must produce infinitely many new states. This is impossible with sets of equation  $E_{\mathcal{F}}^c$  and  $E^r$  as shown in Lemma 4.

## 5.2 Criterion for Functional TRSs

Now, we consider functional programs viewed as TRSs. We assume that such TRSs are left-linear, which is a common assumption on TRSs obtained from functional programs [2]. In this section, we will restrict ourselves to sufficiently complete TRSs obtained from functional programs and will refer to them as *functional TRSs*. For TRSs representing functional programs, defining contracting equations of  $E_{\mathcal{C}}^c$  on  $\mathcal{C}$  rather than on  $\mathcal{F}$  is enough to guarantee termination of completion. This is more convenient and also closer to what is usually done in static analysis where abstractions are usually defined on data and not on function applications. Since the TRSs we consider are sufficiently complete, any term of  $\mathcal{T}(\mathcal{F})$  can be rewritten into a data-term of  $\mathcal{T}(\mathcal{C})$ . As above, using equations of  $E_{\mathcal{C}}^c$  we are going to ensure that the data-terms of the computed languages will be recognized by a bounded set of states. To lift-up this property to  $\mathcal{T}(\mathcal{F})$  it is enough to ensure that  $\forall s, t \in \mathcal{T}(\mathcal{F})$  if  $s \rightarrow_R t$  then  $s$  and  $t$  are recognized by equivalent states. This is the role of the set of equations  $E_R$ .

**Definition 11** ( $E_R$ ). *Let  $\mathcal{R}$  be a TRS, the set of  $\mathcal{R}$ -equations is  $E_R = \{l = r \mid l \rightarrow r \in \mathcal{R}\}$ .*

**Theorem 5.** *Let  $\mathcal{A}_0$  be a tree automaton,  $\mathcal{R}$  a sufficiently complete left-linear TRS and  $E$  a set of equations. If  $E \supseteq E^r \cup E_{\mathcal{C}}^c \cup E_R$  with  $E_{\mathcal{C}}^c$  contracting then completion of  $\mathcal{A}_0$  by  $\mathcal{R}$  and  $E$  terminates.*

*Proof.* Firstly, we can use a reasoning similar to the one used in the proof of Lemma 4 to show that the number of states recognizing terms of  $\mathcal{T}(\mathcal{C})$  are in finite number. Let  $G \subseteq \mathcal{T}(\mathcal{C})$  be the finite set of normal forms of  $\mathcal{T}(\mathcal{C})$  w.r.t.  $\overrightarrow{E_{\mathcal{C}}^c}$ . Since  $E \supseteq E^r \cup E_{\mathcal{C}}^c$ , like in the proof of Lemma 4, we can show that in any completed automaton, terms of  $\mathcal{T}(\mathcal{C})$  are recognized by no more states than terms in  $G$ . Secondly, since  $\mathcal{R}$  is sufficiently complete, for all term  $s \in \mathcal{T}(\mathcal{F}) \setminus \mathcal{T}(\mathcal{C})$  we know that there exists a term  $t \in \mathcal{T}(\mathcal{C})$  such that  $s \rightarrow_{\mathcal{R}}^* t$ . The fact that  $E \supseteq E_R$  guarantees that  $s$  and  $t$  will be recognized by equivalent states in the completed (and simplified) automaton. Since the number of states necessary to recognize  $\mathcal{T}(\mathcal{C})$  is finite, so is the number of states necessary to recognize terms of  $\mathcal{T}(\mathcal{F})$ .

Finally, to exploit the types of the functional program, we now see  $\mathcal{F}$  as a many-sorted signature whose set of sorts is  $\mathcal{S}$ . Each symbol  $f \in \mathcal{F}$  is associated to a profile  $f : S_1 \times \dots \times S_k \mapsto S$  where  $S_1, \dots, S_k, S \in \mathcal{S}$  and  $k$  is the arity of  $f$ . Well-sorted terms are inductively defined as follows:  $f(t_1, \dots, t_k)$  is a well-sorted term of sort  $S$  if  $f : S_1 \times \dots \times S_k \mapsto S$  and  $t_1, \dots, t_k$  are well-sorted terms of sorts  $S_1, \dots, S_k$ , respectively. We denote by  $\mathcal{T}(\mathcal{F}, \mathcal{X})^{\mathcal{S}}$ ,  $\mathcal{T}(\mathcal{F})^{\mathcal{S}}$  and  $\mathcal{T}(\mathcal{C})^{\mathcal{S}}$  the set of well-sorted terms, ground terms and constructor terms, respectively. Note that we have  $\mathcal{T}(\mathcal{F}, \mathcal{X})^{\mathcal{S}} \subseteq \mathcal{T}(\mathcal{F}, \mathcal{X})$ ,  $\mathcal{T}(\mathcal{F})^{\mathcal{S}} \subseteq \mathcal{T}(\mathcal{F})$  and  $\mathcal{T}(\mathcal{C})^{\mathcal{S}} \subseteq \mathcal{T}(\mathcal{C})$ . We assume that  $\mathcal{R}$  and  $E$  are *sort preserving*, i.e. that for all rule  $l \rightarrow r \in \mathcal{R}$  and all equation  $u = v \in E$ ,  $l, r, u, v \in \mathcal{T}(\mathcal{F}, \mathcal{X})^{\mathcal{S}}$ ,  $l$  and  $r$  have the same sort and so do  $u$  and  $v$ . Note that well-typedness of the functional program entails the well-sortedness of  $\mathcal{R}$ . We still assume that the (sorted) TRS is sufficiently complete, which is defined in a similar way except that it holds only for well-sorted terms, i.e. for all  $s \in \mathcal{T}(\mathcal{F})^{\mathcal{S}}$  there exists a term  $t \in \mathcal{T}(\mathcal{C})^{\mathcal{S}}$  such that  $s \rightarrow_{\mathcal{R}}^* t$ . We slightly refine the definition of contracting equations as follows. For all sort  $S$ , if  $S$  has a unique constant symbol we note it  $c^S$ .

**Definition 12** (Set  $E_{\mathcal{K}, \mathcal{S}}^c$  of contracting equations for  $\mathcal{K}$  and  $\mathcal{S}$ ). *Let  $\mathcal{K} \subseteq \mathcal{F}$ . The set of well-sorted equations  $E_{\mathcal{K}, \mathcal{S}}^c$  is contracting (for  $\mathcal{K}$ ) if its equations are of the form (a)  $u = u|_p$  with  $u$  linear and  $p \neq \Lambda$ , or (b)  $u = c^S$  with  $u$  of sort  $S$ , and if the set of normal forms of  $\mathcal{T}(\mathcal{K})^{\mathcal{S}}$  w.r.t. the TRS  $\overrightarrow{E_{\mathcal{K}, \mathcal{S}}^c} = \{u \rightarrow v \mid u = v \in E_{\mathcal{K}, \mathcal{S}}^c \wedge (v = u|_p \vee v = c^S)\}$  is finite.*

The termination theorem for completion of the sorted TRSs is close to the previous one except that it needs  $R/E$ -coherence of  $\mathcal{A}_0$ . This is useful to ensure that terms recognized by completed automata are well-sorted.

**Theorem 6.** *Let  $\mathcal{A}_0$  be a tree automaton recognizing well-sorted terms,  $\mathcal{R}$  a sufficiently complete sort-preserving left-linear TRS and  $E$  a sort-preserving set of equations. If  $E \supseteq E^r \cup E_{\mathcal{C},S}^c \cup E_{\mathcal{R}}$  with  $E_{\mathcal{C},S}^c$  contracting and  $\mathcal{A}_0$  is  $R/E$ -coherent then completion of  $\mathcal{A}_0$  by  $\mathcal{R}$  and  $E$  terminates.*

*Proof.* Let  $\mathcal{A}$  be any tree automaton obtained by completion of  $\mathcal{A}_0$  by  $\mathcal{R}$  and  $E$ . As in Lemma 4, from finiteness of the set normal forms of  $\mathcal{T}(\mathcal{C})^S$  w.r.t.  $\overrightarrow{E_{\mathcal{C},S}^c}$ , we can obtain finiteness of the set of states recognizing terms of  $\mathcal{T}(\mathcal{C})^S$  in the completed automaton. The only slight difference comes from rules of the form  $u = c^S$ . If a term  $s \in \mathcal{T}(\mathcal{C})^S$  is not in normal form w.r.t.  $\overrightarrow{E_{\mathcal{C},S}^c}$  because the rule  $u \rightarrow c^S$  applies then we have:  $s = C[u\sigma] \rightarrow_{\mathcal{A}}^* q$ . Thus there exists a state  $q'$  such that  $u\sigma \rightarrow_{\mathcal{A}}^* q'$ . Since  $c^S$  is the only constant of sort  $S$  and since  $u\sigma$  is of sort  $S$ , we know that  $c^S$  is necessarily a subterm of  $u\sigma$ . Thus there exists a state  $q''$  such that  $c^S \rightarrow_{\mathcal{A}}^* q''$  and since completed automata are simplified,  $q' = q''$  and finally  $C[c^S] \rightarrow_{\mathcal{A}}^* q$ . As in Lemma 4, we can iterate the process until finding a normal form of  $\overrightarrow{E_{\mathcal{C},S}^c}$ . This entails the finiteness of the set of states recognizing terms of  $\mathcal{T}(\mathcal{C})^S$  in  $\mathcal{A}$ . Then, as in the proof of Theorem 5 we can use the fact that  $E \supseteq E_{\mathcal{R}}$  to have that terms of  $\mathcal{T}(\mathcal{F})^S$  are recognized in  $\mathcal{A}$  using a finite set of states. What remains to be proved is that  $\mathcal{A}$  recognizes only well-sorted terms, *i.e.* that it recognizes no term of  $\mathcal{T}(\mathcal{F}) \setminus \mathcal{T}(\mathcal{F})^S$ . This is true because  $\mathcal{A}_0$  is  $R/E$ -coherent, and by Theorem 3,  $\mathcal{L}(\mathcal{A}) \subseteq \mathcal{R}_E^*(\mathcal{L}(\mathcal{A}_0))$ . Besides,  $\mathcal{R}$  and  $E$  being sort-preserving, so is  $\mathcal{R}/E$ . Thus, terms of  $\mathcal{L}(\mathcal{A})$  are all well-sorted.

### 5.3 Experiments

The objective of a data-flow analysis is to predict the set of all program states reachable from a language of initial function calls, *i.e.* to over-approximate  $\mathcal{R}^*(\mathcal{L}(\mathcal{A}))$  where  $\mathcal{R}$  represents the functional program and  $\mathcal{A}$  the language of initial function calls. In this setting, we automatically compute an automaton  $\mathcal{A}_{\mathcal{R},E}^*$  over-approximating  $\mathcal{R}^*(\mathcal{L}(\mathcal{A}))$ . But we can do more. Since we are dealing with left-linear TRS, it is possible to build  $\mathcal{A}_{\text{IRR}(\mathcal{R})}$  recognizing  $\text{IRR}(\mathcal{R})$ . Finally, since tree automata are closed by all boolean operation, we can compute an approximation of all the results of the function calls by computing the tree automaton recognizing the intersection between  $\mathcal{A}_{\mathcal{R},E}^*$  and  $\mathcal{A}_{\text{IRR}(\mathcal{R})}$ .

Here is an example of application of those theorems. Completions are performed using Timbuk. All the  $\mathcal{A}_{\text{IRR}(\mathcal{R})}$  automata and intersections were performed using Tam1. All completion results have been certified by Coq using the Coq-extracted completion checker [7]. All tools are freely available [28]. More details and more examples can be found in [17].

**Ops** append:2 rev:1 nil:0 cons:2 a:0 b:0    **Vars** X Y Z U Xs

**TRS** R

append(nil,X)->X	rev(nil)->nil
append(cons(X,Y),Z)->cons(X,append(Y,Z))	rev(cons(X,Y))->append(rev(Y),cons(X,nil))

**Automaton** A0 **States** q0 qla qlb qnil qf qa qb **Final States** q0 **Transitions**

rev(qla)->q0	cons(qb,qnil)->qlb	cons(qa,qla)->qla	nil->qnil
cons(qa,qlb)->qla	a->qa	cons(qb,qlb)->qlb	b->qb

**Equations** E **Rules**

append(nil,X)=X	a=a   b=b   nil=nil	cons(X,cons(Y,Z))=cons(Y,Z)
append(cons(X,Y),Z)=cons(X,append(Y,Z))	cons(X,Y)=cons(X,Y)	
rev(nil)=nil	append(X,Y)=append(X,Y)	
rev(cons(X,Y))=append(rev(Y),cons(X,nil))	rev(X)=rev(X)	

In this example, the TRS  $\mathcal{R}$  encodes the classical *reverse* and *append* functions. The language recognised by automaton  $\mathcal{A}_0$  is the set of terms of the form  $rev([a, a, \dots, b, b, \dots])$ . Note that there is at least one  $a$  and one  $b$  in the list. We assume that  $\mathcal{S} = \{T, list\}$  and sorts for symbols are the following:  $a : T$ ,  $b : T$ ,  $nil : list$ ,  $cons : T \times list \mapsto list$ ,  $append : list \times list \mapsto list$  and  $rev : list \mapsto list$ . Now, to use Theorem 6, we need to prove each of its assumptions. The set  $E$  of equations contains  $E_{\mathcal{R}}$ ,  $E^r$  and  $E_{\mathcal{C}, \mathcal{S}}^c$ . The set of Equations  $E_{\mathcal{C}, \mathcal{S}}^c$  is contracting because the automaton  $\mathcal{A}_{IRR(\overrightarrow{E_{\mathcal{C}, \mathcal{S}}^c})}$  recognizes a finite language. This automaton can be computed using

Taml: it is the intersection between the automaton  $\mathcal{A}_{\mathcal{T}(C)^S}$ <sup>5</sup> recognising  $\mathcal{T}(C)^S$  and the automaton  $\mathcal{A}_{IRR(\{cons(X, cons(Y, Z)) \rightarrow cons(Y, Z)\})}$ :

**States** q2 q1 q0 **Final States** q0 q1 q2 **Transitions** b->q2 a->q2 nil->q1 cons(q2, q1)->q0

The language of  $\mathcal{A}_0$  is well-sorted and  $E$  and  $\mathcal{R}$  are sort preserving. We can prove sufficient completeness of  $\mathcal{R}$  on  $\mathcal{T}(\mathcal{F})^S$  using, for instance, Maude [8]. The last assumption of Theorem 6 to prove is that  $\mathcal{A}_0$  is  $R/E$ -coherent. This can be shown by remarking that each state  $\mathcal{A}_0$  recognizes at least one term and that for all state  $q$  such that  $s \rightarrow_{\mathcal{A}_0}^{\ell^*} q$  and  $t \rightarrow_{\mathcal{A}_0}^{\ell^*} q$  then  $s =_E t$ . For instance  $cons(b, cons(b, nil)) \rightarrow_{\mathcal{A}_0}^{\ell^*} qlb$  and  $cons(b, nil) \rightarrow_{\mathcal{A}_0}^{\ell^*} qlb$  and  $cons(b, cons(b, nil)) =_E cons(b, nil)$ . Thus, completion is guaranteed to terminate: after 4 completion steps (7 ms) we obtain a fixpoint automaton  $\mathcal{A}_{\mathcal{R}, E}^*$  with 11 transitions. To restrain its language to normal forms it is necessary to compute the intersection with  $IRR(R)$ . Since we are dealing with sufficiently complete TRSs, we know that  $IRR(R) \subseteq \mathcal{T}(C)^S$ . Thus, we can use again  $\mathcal{A}_{\mathcal{T}(C)^S}$  for the intersection that is:

**States** q3 q2 q1 q0 **Final States** q3 **Transitions**  
a->q0 nil->q1 b->q2 cons(q0, q1)->q3 cons(q0, q3)->q3 cons(q2, q1)->q3 cons(q2, q3)->q3

which recognizes any (non empty) flat list of  $a$  and  $b$ . Thus, our analysis preserved the property that the result cannot be the empty list but lost the order of the elements in the list. This is not surprising because the equation  $cons(X, cons(Y, Z)) = cons(X, Z)$  makes  $cons(a, cons(b, nil))$  equal to  $cons(a, nil)$ . It is possible to refine by hand  $E_{\mathcal{C}, \mathcal{S}}^c$  as follows:

$cons(a, cons(a, X)) = cons(a, X)$ ,  $cons(b, cons(b, X)) = cons(b, X)$ ,  $cons(a, cons(b, cons(a, X))) = cons(a, X)$

This set of equations avoids the previous problem. Again,  $E$  verifies the conditions of Theorem 6 and completion is still guaranteed to terminate. The result is the automaton  $\mathcal{A}_{\mathcal{R}, E}^{t*}$  having 19 transitions. This time, intersection with  $\mathcal{A}_{\mathcal{T}(C)^S}$  gives:

**States** q4 q3 q2 q1 q0 **Final States** q4 **Transitions**  
a->q1 b->q3 nil->q0 cons(q1, q0)->q2 cons(q1, q2)->q2 cons(q3, q2)->q4 cons(q3, q4)->q4

This automaton exactly recognizes lists of the form  $[b, b, \dots, a, a, \dots]$  with at least one  $b$  and one  $a$ , as expected. Hopefully, refinement of equational approximations can be automatized [3] and can be used in this setting, see [17] for examples.

## 6 Conclusion and further research

In this paper we defined a criterion on the set of approximation equations to guarantee termination of tree automata completion. When dealing with, so called, functional TRS this criterion is close to what is generally expected from an abstract domain used for static analysis: define a finite model for an infinite set of data-terms. This is a first step to use reachability analysis techniques of rewriting for static analysis of functional programs. For this technique to be applicable on real programs there remains some interesting points to address.

<sup>5</sup> Such an automaton has one state per sort and one transitions per constructor. For instance, on our example  $\mathcal{A}_{\mathcal{T}(C)^S}$  will have transitions:  $a \rightarrow qT$ ,  $b \rightarrow qT$ ,  $cons(qT, qlist) \rightarrow qlist$  and  $nil \rightarrow qlist$ .

*Dealing with higher-order functions.* Higher-order functions can be encoded into first order TRS using a simple encoding borrowed from [20]: defined symbols become constants, constructor symbols remain the same, and an additional *application* operator '@' of arity 2 is introduced. On all the examples of [25], which is the state of the art of data-flow analysis of higher-order functional programs, it has been shown that using this simple encoding completion produces exactly the same results [17].

*Dealing with evaluation strategies.* The technique proposed here, as well as [25], over-approximate the set of results for all evaluation strategies. As far as we know, no static analysis technique for functional programs can take into account evaluation strategies. Thus, if a program is not terminating using the usual call-by-value strategy (innermost rewriting strategy) but terminating by call-by-need (outermost rewriting strategy plus sharing) the analysis will give the results obtained by call-by-need. However, it has been shown that it is possible to restrict the completion algorithm to recognize only innermost descendants [16], *i.e.* call-by-value results. If the approximation is precise enough, any non terminating program with call by value will have an empty set of results. An interesting open research direction is to build from those results a criterion for non termination of functional programs by call-by-value.

*Dealing with built-in types.* Values manipulated by *real* functional programs are not always terms or trees. They can be numerals or be terms embedding numerals. In [14], it has been shown that completion can compute over-approximations of reachable terms embedding built-in terms. The structural part of the term is approximated using tree automata and the built-in part is approximated using lattices and abstract interpretation.

*Presenting the results of the analysis.* Our long term objective is to define a static analysis tool complementary to the usual type inference tools used by modern functional programs compilers. The computed tree automata can give an information that is more precise than a type. For instance, it can discriminate between an empty and a non empty list. An open question is how to present the computed information, *i.e.* a tree automaton, to the user so that he can find a potential problem in the function he has defined.

Besides, there remain some interesting theoretical points to solve. In section 5, we saw that having a set of equations such that  $\mathcal{T}(\mathcal{F})/_{{=}_E}$  is finite is not enough to guarantee termination of completion. This is due to the fact that the simplification algorithm does not merge all states recognizing *E*-equivalent terms. Having a simplification algorithm ensuring this property is not trivial. First, the theory defined by *E* has to be decidable. Second, since it is not possible, in general, to finitely enumerate all the terms recognized by all the states of a tree automaton, how to find all the *E*-equivalent terms recognized by the automaton? This is an open problem.

Similarly, proving that  $\mathcal{T}(\mathcal{F})/_{{=}_E}$  is finite is an open problem. This question is not decidable in general [29]. However, defining simple criteria on *E* for  $\mathcal{T}(\mathcal{F})/_{{=}_E}$  to be finite is also an open interesting problem. For instance, if *E* can be oriented into a TRS  $\mathcal{R}$  which is terminating, confluent and such that  $\text{IRR}(\mathcal{R})$  is finite then  $\mathcal{T}(\mathcal{F})/_{{=}_E}$  is finite [29].

## 7 Conclusion

## References

1. A. Armando, D. Basin, Y. Boichut, Y. Chevalier, L. Compagna, J. Cuellar, P. Hankes Drielsma, P.-C. Héam, O. Kouchnarenko, J. Mantovani, S. Mödersheim, D. von Oheimb, M. Rusinowitch, J. Santos Santiago, M. Turuani, L. Viganò, and L. Vigneron. The AVISPA Tool for the automated validation of internet security protocols and applications. In *CAV'2005*, volume 3576 of *LNCS*, pages 281–285. Springer, 2005.
2. F. Baader and T. Nipkow. *Term Rewriting and All That*. Cambridge University Press, 1998.

3. Y. Boichut, B. Boyer, T. Genet, and A. Legay. Equational Abstraction Refinement for Certified Tree Regular Model Checking. In *ICFEM'12*, volume 7635 of *LNCS*. Springer, 2012.
4. Y. Boichut, R. Courbis, P.-C. Héam, and O. Kouchnarenko. Handling non left-linear rules when completing tree automata. *IJFCS*, 20(5), 2009.
5. Y. Boichut, T. Genet, T. Jensen, and L. Leroux. Rewriting Approximations for Fast Prototyping of Static Analyzers. In *RTA*, volume 4533 of *LNCS*, pages 48–62. Springer, 2007.
6. Y. Boichut, P.-C. Héam, and O. Kouchnarenko. Automatic Approximation for the Verification of Cryptographic Protocols. In *Proc. AVIS'2004, joint to ETAPS'04, Barcelona (Spain)*, 2004.
7. B. Boyer, T. Genet, and T. Jensen. Certifying a Tree Automata Completion Checker. In *IJCAR'08*, volume 5195 of *LNCS*. Springer, 2008.
8. Manuel Clavel, Francisco Durán, Steven Eker, Patrick Lincoln, Narciso Martí-Oliet, José Meseguer, and José F. Quesada. Maude homepage, 2009. <http://maude.cs.uiuc.edu>.
9. H. Comon, M. Dauchet, R. Gilleron, F. Jacquemard, D. Lugiez, C. Löding, S. Tison, and M. Tommasi. Tree automata techniques and applications. <http://tata.gforge.inria.fr>, 2008.
10. G. Feuillade, T. Genet, and V. Viet Triem Tong. Reachability Analysis over Term Rewriting Systems. *Journal of Automated Reasoning*, 33 (3-4):341–383, 2004.
11. T. Genet. Decidable approximations of sets of descendants and sets of normal forms. In *Proc. 9th RTA Conf., Tsukuba (Japan)*, volume 1379 of *LNCS*, pages 151–165. Springer-Verlag, 1998.
12. T. Genet. Reachability analysis of rewriting for software verification. Université de Rennes 1, 2009. Habilitation document, <http://www.irisa.fr/celtique/genet/publications.html>.
13. T. Genet and F. Klay. Rewriting for Cryptographic Protocol Verification. In *Proc. 17th CADE Conf., Pittsburgh (Pen., USA)*, volume 1831 of *LNAI*. Springer-Verlag, 2000.
14. T. Genet, T. Le Gall, A. Legay, and V. Murat. A Completion Algorithm for Lattice Tree Automata. In *CIAA'13*, volume 7982 of *LNCS*, pages 134–145, 2013.
15. T. Genet and R. Rusu. Equational tree automata completion. *Journal of Symbolic Computation*, 45:574–597, 2010.
16. T. Genet and Y. Salmon. Proving reachability properties on term rewriting systems with strategies. In *2nd Joint International Workshop on Strategies in Rewriting, Proving and Programming, IWS'12*, London, 2012.
17. T. Genet and Y. Salmon. Tree Automata Completion for Static Analysis of Functional Programs. Technical report, INRIA, 2013. <http://hal.archives-ouvertes.fr/hal-00780124/PDF/main.pdf>.
18. A. Geser, D. Hofbauer, J. Waldmann, and H. Zantema. On tree automata that certify termination of left-linear term rewriting systems. In *RTA'05*, volume 3467 of *LNCS*, pages 353–367. Springer, 2005.
19. F. Jacquemard. Decidable approximations of term rewriting systems. In H. Ganzinger, editor, *Proc. 7th RTA Conf., New Brunswick (New Jersey, USA)*, pages 362–376. Springer-Verlag, 1996.
20. N. D. Jones. Flow analysis of lazy higher-order functional programs. In S. Abramsky and C. Hankin, editors, *Abstract Interpretation of Declarative Languages*, pages 103–122. Ellis Horwood, Chichester, England, 1987.
21. Naoki Kobayashi. Model checking higher-order programs. *J. ACM*, 60(3):20, 2013.
22. J. Kochems and L. Ong. Improved Functional Flow and Reachability Analyses Using Indexed Linear Tree Grammars. In *RTA'11*, volume 10 of *LIPICs*. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2011.
23. A. Lisitsa. Finite Models vs Tree Automata in Safety Verification. In *RTA'12*, volume 15 of *LIPICs*, pages 225–239, 2012.
24. A. Middeldorp. Approximations for strategies and termination. *ENTCS*, 70(6):1–20, 2002.
25. L. Ong and S. Ramsay. Verifying higher-order functional programs with pattern-matching algebraic data types. In *POPL'11*, 2011.
26. T. Takai. A Verification Technique Using Term Rewriting Systems and Abstract Interpretation. In *Proc. 15th RTA Conf., Aachen (Germany)*, volume 3091 of *LNCS*, pages 119–133. Springer, 2004.
27. T. Takai, Y. Kaji, and H. Seki. Right-linear finite-path overlapping term rewriting systems effectively preserve recognizability. In *Proc. 11th RTA Conf., Norwich (UK)*, volume 1833 of *LNCS*. Springer-Verlag, 2000.
28. Timbuk – reachability analysis and Tree Automata Calculations. IRISA / Université de Rennes 1, 2012. <http://www.irisa.fr/celtique/genet/timbuk/>.
29. S. Tison. Finiteness of the set of  $E$ -equivalence classes is undecidable, 2010. Private communication.